

-1-

Date: 11/28/00 Express Mail Label No. EL552571276US

Inventor(s): Simon Kasif, Beth T. Logan, Pedro J. Moreno and Baris E. Suzek

Attorney's Docket No.: 0918.2033-000 (P00-3373)

COMPUTER METHOD AND APPARATUS FOR UNIFORM
REPRESENTATION OF GENOME SEQUENCES

BACKGROUND OF THE INVENTION

Computational methods for biological sequence analysis are playing an increasingly important role in biology and medicine. The key question addressed by these methods is the discovery of the function of a protein or gene. It is well known that the function of a protein is dictated by its amino acid sequence since this determines the structure of the protein and thus its interaction with the environment.

Proteins are the building blocks of life, supporting a variety of functions which are essential for cell life. These include protection from infections or cancers, gene regulation, survival in different conditions, growth, differentiation, regeneration and others. In fact, the function of every cell in a living organism (whether microbial or human) is determined by which proteins (genes) are expressed in the cell and how they interact in the particular cell environment.

15 The area of protein function is particularly timely because the new technology of high-throughput genomics generates thousands of hypothetical genes that have not been assigned a putative function. There are numerous commercial applications. Classifying new genes into categories opens many opportunities for new medical treatments. Genes are often used as drugs directly (e.g., insulin), or drug targets (e.g., attacking a particular

gene in a microbial organism). Other applications include the design of pesticides, design of new crops, gene therapies and rational drug design.

Proteins are macromolecules found in living organisms which play many roles essential to sustaining life (e.g., forming the physical framework of the organism, acting as enzymes to promote chemical reactions). A protein is composed of a sequence of several hundred amino acids. Proteins are created in living cells by translating the coding regions (genes) of the DNA sequence. Different proteins are expressed in different cells. The level of expression of different cells determines the cell function. Since proteins are long and linear complex molecules, they "fold" to give a 3D shape.

Biologists have identified four levels of structure which can influence the protein's function:

1. Primary structure--the sequence of amino acids
2. Secondary structure--the presence or absence of small "sub-folds".
These are regular patterns formed by local folding of the protein (e.g., helices and sheets).
3. Tertiary structure--the final 3D shape
4. Quaternary structure--complexes formed with other proteins.

Given one level of structure, it is not necessarily a trivial task to predict the next level. Hence, function prediction from the primary structure alone is difficult.

Therefore, techniques other than sequencing are needed to determine the 3D structure and ultimately the protein function.

The traditional and still most reliable way to perform protein structure prediction is to use laboratory-based techniques such as X-ray crystallography. However, recent years have seen the development of software-based solutions. One such technique is to use dynamic programming-based alignment tools such as "BLAST" to match the new sequence to previously labeled protein sequences (Altshul et al., 1990, Basic Local Alignment Search Tool, JMB 215:403-410). Alternatively, statistical techniques such as Hidden Markov Models (HMM's) can be used to build a model for each labeled class (E. Sonnhammer, S. Eddy and R. Durbin, "Pfam: A Comprehensive Database of Protein

0918.2033-000

Sub
a.

Families Based on Seed Alignments," *Proteins*, 1997, pages 405-420). (A. Krogh, M. Brown, I. Mian, K. Sjolander and D. Haussler, "Hidden Markov Models in Computational Biology: Applications to Protein Modeling, *J. of Molecular Biology*, 1994, Volume 235, 1501-1531.) Still another alternative is to learn the boundaries
5 between protein classes rather than a model for the class itself. (Jaakkola, Diekhans, Haussler, "Using the Fisher kernel method to detect remote protein homologies," in Proceedings of ISBM '99). The first two approaches use the protein sequence itself directly to perform classification. The last one uses a HMM to compute the gradient of the protein being produced by the HMM with respect to each of the parameters of the
10 HMM. In summary, none of these methods uses the sensitivity of parts of the protein to motifs to build a feature vector.

Lab-based techniques, such as X-ray crystallography, are expensive and time-consuming. In addition, X-ray crystallography relies on having relatively large amounts of the protein. It cannot work with just a primary description of the protein (i.e., the
15 sequence of amino acids in a file). Finally, it is not possible to crystallize certain proteins in any case (e.g., membrane spanning proteins).

BLAST and other dynamic programming methods are more time-consuming and less accurate than statistical-based techniques.

SUMMARY OF THE INVENTION

20 The invention addresses the problem of classifying, clustering or indexing proteins and other biological sequences such as genes by using an alternative representation based on high dimensional vectors. Each of the components of the vector represents the sensitivity of the protein (or sequence) to a particular biological motif (described later). Once obtained, this new representation can be used in conjunction
25 with many existing machine learning techniques to analyze the sequences of interest. For example, this new representation may be combined with discriminative classification methods to classify new proteins from the amino acid sequence alone.

The following discloses a new representation of proteins (genes) as objects in a very high-dimensional vector space. This representation offers numerous opportunities for predictive analysis of the space of biological sequences in a novel fashion deploying high-dimensional analysis techniques. The representation relies on aligning very short
5 motif elements (biological templates) to the protein sequence. Subsequently, each protein is encoded as a multi-dimensional vector X , where dimension X_i corresponds to the score obtained by obtaining the maximum score of scoring (convolving) element E_i "against" the protein. The representation allows the use of existing templates (motifs) or to "train" new ones.

10 For example, currently, limited databases exist which contain protein domain sequences (primary structure) annotated with their secondary and tertiary structure. A protein domain is a subsequence of interest found in proteins. One use of the present invention is to use this labeled data to build models for known protein structures, and then to automatically annotate new proteins according to the models. However, the
15 general idea of the invention may also apply to other protein or gene classification problems and to cluster or index biological sequences.

In a preferred embodiment, a method and apparatus transforms typically differing length text string representations (i.e., sequences) of biological fragments into uniform length representations. A comparison database stores a predefined number of
20 known biological sequences. A comparison routine compares and scores a subject sequence against each known sequence in the database. Each individual score (one for each known sequence in the database) serves as a vector element forming a fixed length vector representation of the subject sequence. Vector length equals the predefined number of known biological sequences in the database. Scoring is by a counting of the
25 number of times the known biological sequence is found in the subject sequence, or the probability of the subject sequence being generated by the known biological sequence.

BRIEF DESCRIPTION OF THE DRAWINGS

The foregoing and other objects, features and advantages of the invention will be apparent from the following more particular description of preferred embodiments of the invention, as illustrated in the accompanying drawings in which like reference
5 characters refer to the same parts throughout the different views. The drawings are not necessarily to scale, emphasis instead being placed upon illustrating the principles of the invention.

Fig. 1 is a block diagram of a computer implementation of the present invention.

Fig. 2 is a flow diagram of the present invention overall process.

10 Fig. 3 is a schematic of the invention transform into feature vectors utilized in the overall process of Fig. 2.

DETAILED DESCRIPTION OF THE INVENTION

By way of overview, a cell has an operational center called the nucleus which contains structures called chromosomes. Chemically, chromosomes are formed of
15 deoxyribonucleic acid (DNA) and associated protein molecules. Structurally, each chromosome has tens of thousands of genes. Some genes are referred to as "encoding" (or carrying information for constructing) proteins which are essential in the structuring, functioning and regulating of cells, tissues and organs. Thus, for each organism, the components of the DNA molecules encode all the information necessary for creating
20 and maintaining life of the organism. See Human Genome Program, U.S. Department of Energy, "Primer on Molecular Genetics", Washington, D.C., 1992.

The shape of a DNA molecule can be thought of as a twisted ladder. That is, the DNA molecule is formed of two parallel side strands of sugar and phosphate molecules connected by orthogonal/cross pieces (rungs) of nitrogen-containing chemicals called
25 bases. Each long side strand is formed of a particular series of units called nucleotides. Each nucleotide comprises one sugar, one phosphate and a nitrogenous base. The order of the bases in this series (the side strands series of nucleotides) is called the DNA sequence.

5

10

15

20

25

Different amino acid sequences have different length text string representations.

predefined set 17 of known biological fragments, the invention software program 15

The comparison routine 19 effectively transforms the traditional ACTG text.

is, the comparison routine 19 transforms the input sequences of varying length into

In the preferred embodiment, the number of known biological fragments in the predefined set 17 defines the length of resulting feature vectors 23.

The output 21 of the invention software 15 (i.e., normalized representations of amino acid sequences, each representation being of the same length) may then be fed
15 into analyses of typical interest in biotechnology. Such analyses include classification, clustering and indexing.

It is understood that input amino acid sequences 11 may be received from input devices (e.g., a keyboard, mouse, etc.), another computer coupled across a communication channel to digital processor 13 (i.e., in a local area, wide area and/or global/Internet network), and the like. Similarly, output 21 of the uniform length feature vectors 23 of the invention software 15 may be transmitted to a data file/data store, another program/processor routine, another computer coupled across a communication channel to digital processor 13, and the like.

Accordingly, the present invention method provides a two-step process 39.

25 First, the invention method converts the amino acid sequences 11 of interest to high dimensional feature vectors 23. Once this transformation has taken place, then one may apply any number of statistical learning techniques to train models for classification, clustering or indexing the protein sequences in the second step of the overall invention process 39. Figs. 2 and 3 describe these steps as detailed below. Although this

description, details the overall process 39 as it applies to the analysis of protein sequences or subsequences, it is understood that invention method and techniques may also be applied to DNA sequences or subsequences.

The first half/phase of the invention method/process 39 illustrated in Fig. 2
 5 converts each protein sequence or subsequence of interest 11 to a new representation of fixed length, i.e., any protein sequence no matter how long it is, is converted into a feature vector 23 of fixed length. Preferably each dimension of these feature vectors 23 represents the sensitivity of the protein to a particular biological motif. Therefore, in order to create feature vectors 23, the invention method first creates or obtains a
 10 comparison database 17 of short, highly conserved regions in related protein domains (step 31). Such regions are often called "blocks", "motifs" or "probabilistic templates".

A working motif is preferably represented by a K by L matrix M in which each of the K rows represents a particular amino acid (or nucleotide for DNA sequences) and L represents the length of the motif. For protein sequences, K = 20. For DNA
 15 sequences K = 4. Each cell, as indicated by [amino acid, position in the length], in the matrix M holds a value that represents the probability of that amino acid existing in that position. This matrix may alternatively store log-ratios rather than probabilities. Thus, a motif may be thought of as a 0-th order Markov model.

The BLOCKS database (Steven Henikoff and Jorja G. Henikoff, "Automated
 20 assembly of protein blocks for database searching," *Nucleic Acids Research*, 19:23, pp. 6565-6572 (1991)) is an example of a database 17 of motifs. Emitof
 (http://dna.stanford.edu/emotif/), and PRINTs (http://bioinf.man.ac.uk.dbbrowser/PRINTS/) are other such databases. These and other published databases may be used as the working predefined set/comparison database 17 in the present invention.
 25 Alternatively, it is possible to create a new motif database 17 from any protein database which has been labeled according to some parameter (e.g., structure). This is achieved by using multiple alignment software to find short multiply aligned ungapped sequences and then collecting statistics about these in a matrix (http://www2.ebi.ac.uk/clustalw/, http://www.blocks.fhcrc.org/). By creating a motif database 17 specific to the proteins

SUB 31

of interest 11, more meaningful feature vectors 23 may be obtained since the motifs from a more general database may not occur in the proteins of interest.

To create a feature vector 23 for each protein sequence 11 of interest, the invention method at step 33 searches for each motif (generated in step 31 and stored in database 17) in the sequence 11 and scores the search results as a count of number of matches found or as a probability, or the like. In the preferred embodiment, in step 33, each motif of length L is scored against the subject protein sequence 11 by computing the probability of every subsequence of length L in the subject sequence 11 being generated by the model (matrix M discussed above) that corresponds to the motif.

This is illustrated in Fig. 3 where subject protein sequence 11 is shown being scored against each motif in comparison database 17 (obtained from step 31). The score (probability or count, etc.) 29 of a first motif against input sequence 11 is indicated as B_1 in Fig. 3. The score of a second motif relative to the same input sequence 11 is indicated as B_2 and so on in Fig. 3. The ordered series of individual motif scores B_i is $[B_1 \dots B_N]$ and represents the feature vector 23 created for subject sequence 11. N is the fixed number of motifs in comparison database 17 that are processed against each input sequence 11 of interest.

Thus, the result at 35 in Fig. 2 is an N -dimensional feature vector where N is the total number of motifs in comparison database 17 as explained above. Each dimension J contains a score describing the degree of alignment of motif J to the subject input sequence 11. For the case where a motif is detected multiple times in input subject sequences 11, the preferred embodiment applies a variety of heuristics at step 35. For example, the invention process 39 takes the maximum of all scores for that block in an input subject sequence 11 or the sum of such scores. In preliminary experiments, Applicants found that taking the maximum score gives superior classification performance. Invention process 39 may also apply a threshold such that scores below a certain number are set to zero at step 35. Additionally, given the complete set of feature vectors 23 for input subject sequences 11, one may (at step 35) reduce the

dimensionality of these vectors using standard dimension reduction techniques such as Principal Components Analysis (PCA).

Continuing in Fig. 2, the second phase in invention overall process 39 includes clustering 34, classification 37 and indexing 30 analyses of interest.

- 5 Once all the protein sequences or subsequences of interest 11 have been transformed to feature vectors 23, models may be generated to describe these features and perform clustering 34, classification 37 or indexing 39. Each of these analyses is described below.

Clustering 34

- 10 A clustering process 34 groups together proteins (subject sequences) 11 with similar feature vectors 23 in order to discover previously unknown relationships between them. For example, using well known algorithms such as k-means or nearest neighbors, it is possible to decide if two proteins 11 as represented by the newly generated feature vectors 23 are close in sequence pattern or not. The key concept here
- 15 is that the new representation (uniform length feature vector 23) allows subsequent analyses to compare proteins (sequences) both reliably and effectively.

Classification 37

- The process of classification 37 attempts to learn a relationship or model given a set of labeled feature vectors 23 called the "training set". Each label denotes the class
- 20 that the vector 23 belongs to. For example, the classes may be defined by protein structural information. Possibly the labeling is generated by clustering. Given this model, unseen vectors, usually denoted the "testing set", are assigned labels according to the models learned. An example of the classification of proteins into structural classes is described below.

Example

1. Given a set of training protein sequences labeled according to structure, convert each of these into a multidimensional feature vector 23 as described above. Utilize the BLOCK's motif database as the comparison database 17 to create the feature vectors 23.
2. Given the labeled feature vectors generated in step 1, learn corresponding Support Vector Machine (SVM) classifiers (Burger, 1998, "A tutorial on Support Vector Machines for Pattern Recognition," *Data Mining and Knowledge Discovery Journal*) to separate each structural class from "the rest of the world". A SVM classifier learns a separating hyperplane between two classes which maximizes the "margin"--the distance between the hyperplane and the nearest datapoint of each class.

The appeal of SVM's is twofold. First, they do not require any complex tuning of parameters, and second they exhibit a great ability to generalize given a small

training corpra. They are particularly amenable for learning in high dimensional spaces. The only parameters needed to tune a SVM are the "capacity" and the choice of kernel. The capacity allows one to control how much tolerance for errors in the classification of training samples one allows and therefore the generalization ability of the SVM. A

5 SVM with high capacity will classify all training samples correctly but will not be able to generalize well for testing samples. In effect, it will construct a classifier too tuned for the training samples which will limit its ability to generalize later on when testing samples are presented to the system. Conversely, a very low capacity will produce a classifier that does not fit the data sufficiently accurately. It will allow many training

10 and testing samples to be classified incorrectly.

The second tuning parameter, called the kernel, allows the SVM to create hyperplanes in high dimensional spaces that effectively separate the training data. Often in the input space training vectors cannot be separated by a simple hyperplane. The kernel allows transforming the data from one space to another space where a simple

15 hyperplane can effectively separate the data in two classes.

In step 2, tune these two parameters separately for each structural family of interest.

An additional step consists of tuning the operating point of the classifier so that one may control the amount of false negatives. In one implementation, Applicants find

20 a threshold value such that any score returned by the SVM that is bigger than this guarantees no false negatives.

3. Given a set of unlabeled structural sequences (the input testing set) convert each of these into a corresponding multidimensional feature vector 23 using BLOCKS as above.

25 4. Now, for each unlabeled feature vector, to determine if it belongs to a particular class, test it using the SVM created for that class. The SVM classifier will produce a "score" representing the distance of the testing feature vector from the margin.

The bigger the score the further away the vector is from the margin and the more confident the classifier is in its own output. If the score is below the threshold set in Step 2, classify the vector (and hence the corresponding test input sequence) as belonging to that particular class. Otherwise, it is classified as not belonging to the
5 class. For multi-class classification one may use standing procedures such as classifying based on the highest score returned by each of the individual classifiers.

While this invention has been particularly shown and described with references to preferred embodiments thereof, it will be understood by those skilled in the art that various changes in form and details may be made therein without departing from the
10 scope of the invention encompassed by the appended claims.

For example, the foregoing describes a method and apparatus for transforming representations of protein or DNA sequences and/or subsequences. It is understood that representations of other biological sequences (human or other) may similarly be transformed using the disclosed techniques and methods.

0918.2033-000